

MODERN PROBABILISTIC MODELING FOR MASSIVE DATA

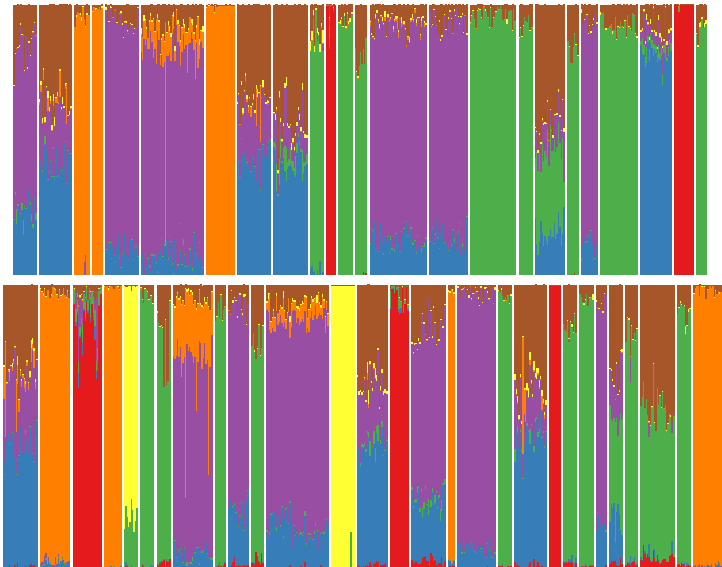
David M. Blei
Columbia University

Modern probabilistic modeling:

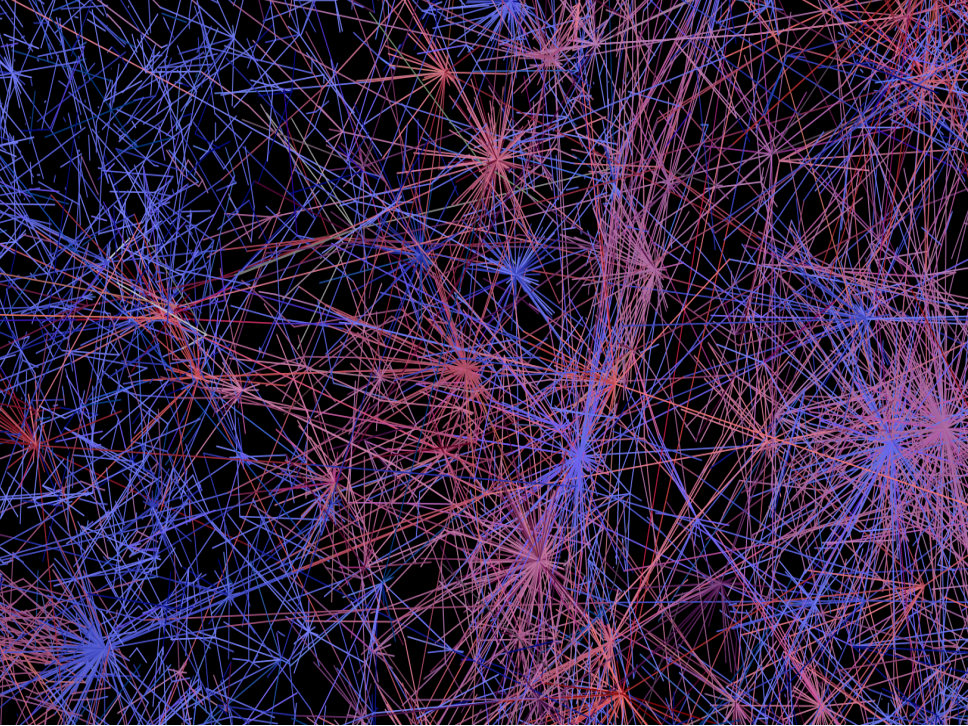
An efficient framework for discovering useful patterns in massive data.

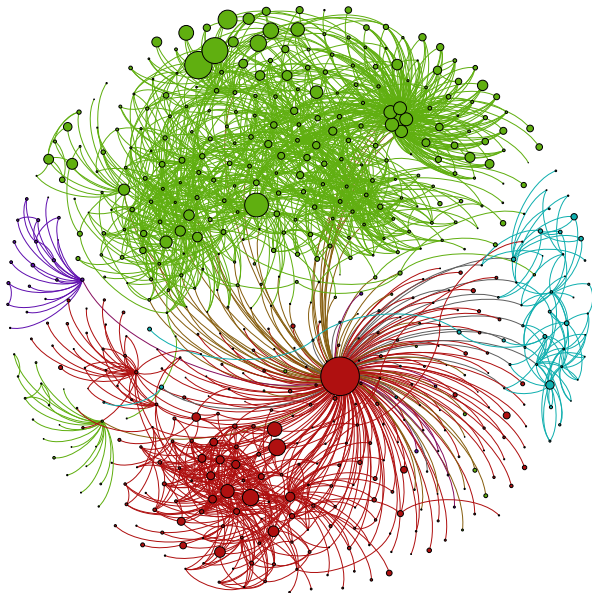


[Amy Pettingill]



Population analysis of 2 billion genetic measurements



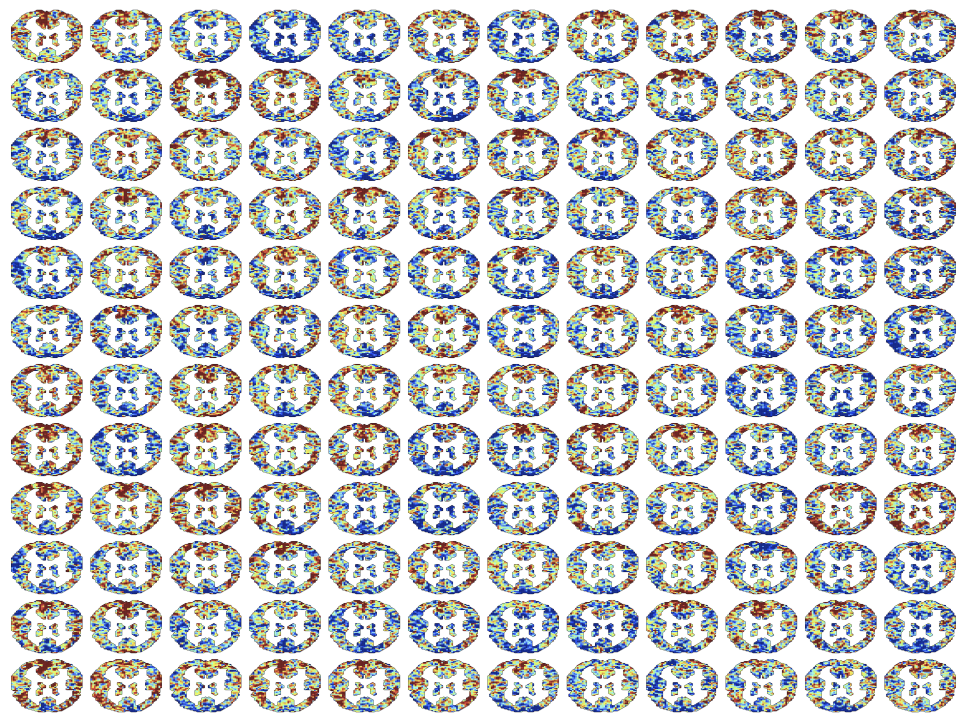


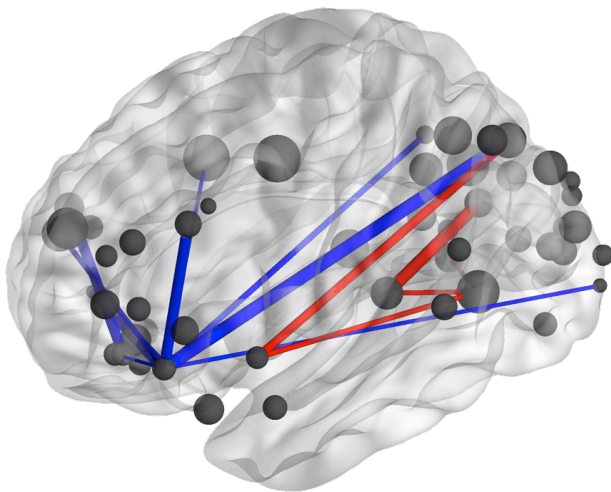
Communities discovered in a 3.7M node network of U.S. Patents





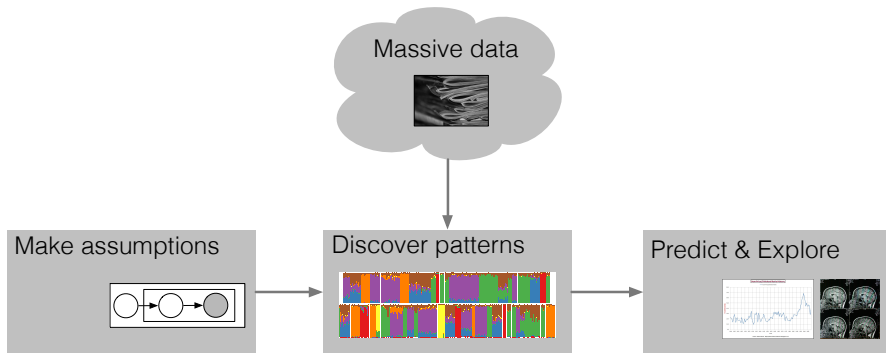
Topics found in 1.8M articles from the New York Times





Neuroscience analysis of 220 million fMRI measurements

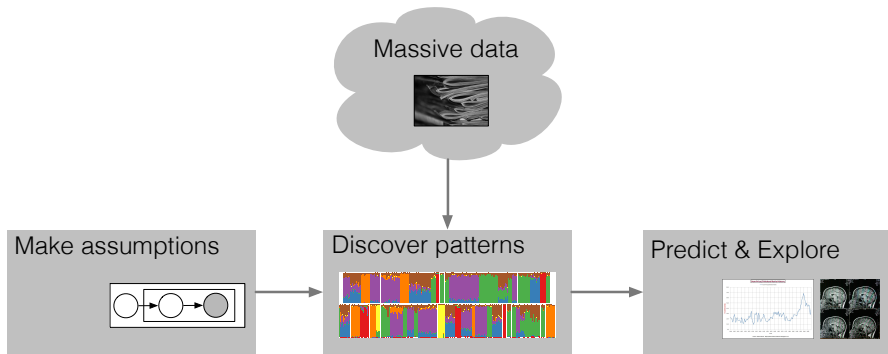
The probabilistic modeling pipeline



Our perspective:

- ▶ This is a framework for **customized data analysis**, crucial to many fields.
- ▶ The pipeline separates assumptions, computation, application
- ▶ It facilitates solving modern data science problems.

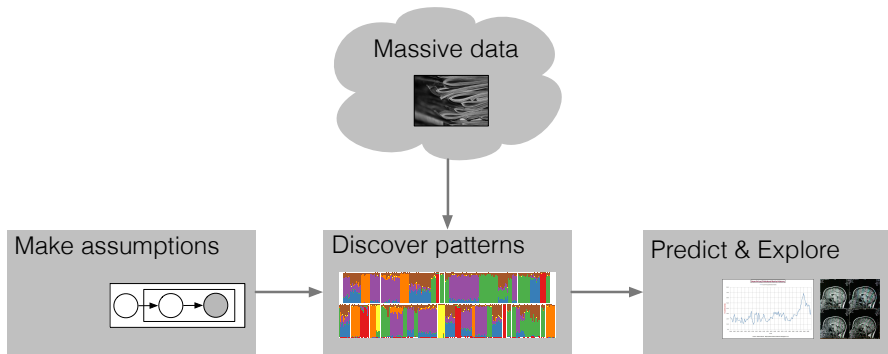
The probabilistic modeling pipeline



Our goal:

- Develop modeling into a **flexible**, **powerful** and **easy-to-use** way to solve real-world problems.

The probabilistic modeling pipeline



Our challenges:

- ▶ Develop new ways to build **flexible models**
- ▶ Develop algorithms that work on many problems and with **massive data**.
- ▶ Solve **new problems** in science, industry, and government



Jaan Aaltosar



Allison Chaney



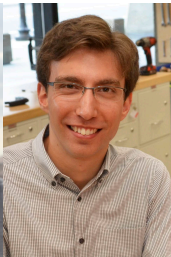
Rajesh Ranganath



Maja Rudolph



Laurent Charlin



Alp Kucukelbir



Stephan Mandt



Jeremy Manning



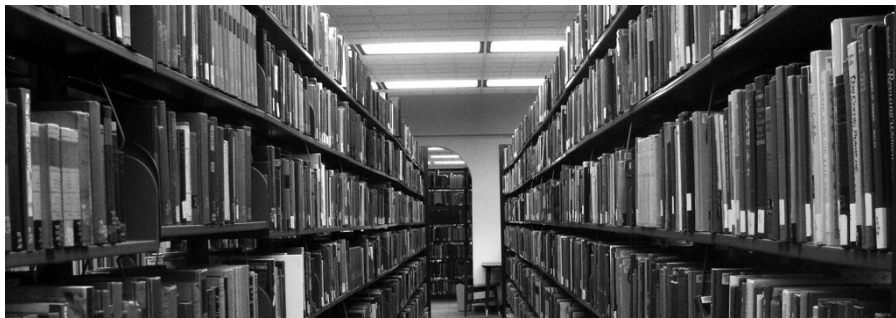
James McInerney

Probabilistic topic models

Powerful and flexible algorithms for analyzing massive collections of text

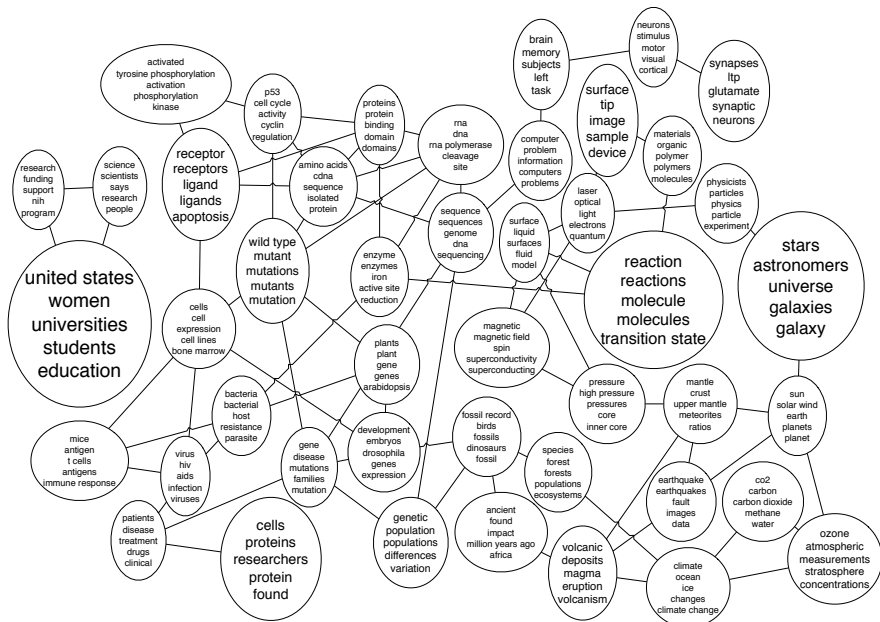


- ▶ **ORGANIZE**
- ▶ **VISUALIZE**
- ▶ **SUMMARIZE**
- ▶ **SEARCH**
- ▶ **PREDICT**
- ▶ **UNDERSTAND**



TOPIC MODELING

1. **Discover** the thematic structure
2. **Annotate** the documents
3. **Use** the annotations to visualize, organize, summarize, ...



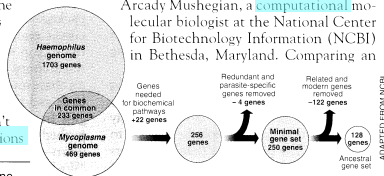
Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic numbers game**, particularly as more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains

Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



Stripping down. **Computer analysis** yields an estimate of the minimum modern and ancient genomes.

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Documents exhibit multiple topics.

Topics

gene 0.04
dna 0.02
genetic 0.01
...

life 0.02
evolve 0.01
organism 0.01
...

brain 0.04
neuron 0.02
nerve 0.01
...

data 0.02
number 0.02
computer 0.01
...

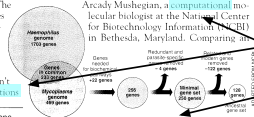
Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here, two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson, a geneticist at the University of Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic** numbers game, particularly if more and more **genomes** are completely mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

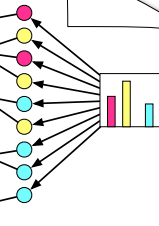


* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



Latent Dirichlet Allocation

Topics



Documents

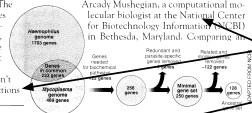
Topic proportions and assignments

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes to the human genome, notes Siv Anderson, a biologist at the University of Warwick, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing the



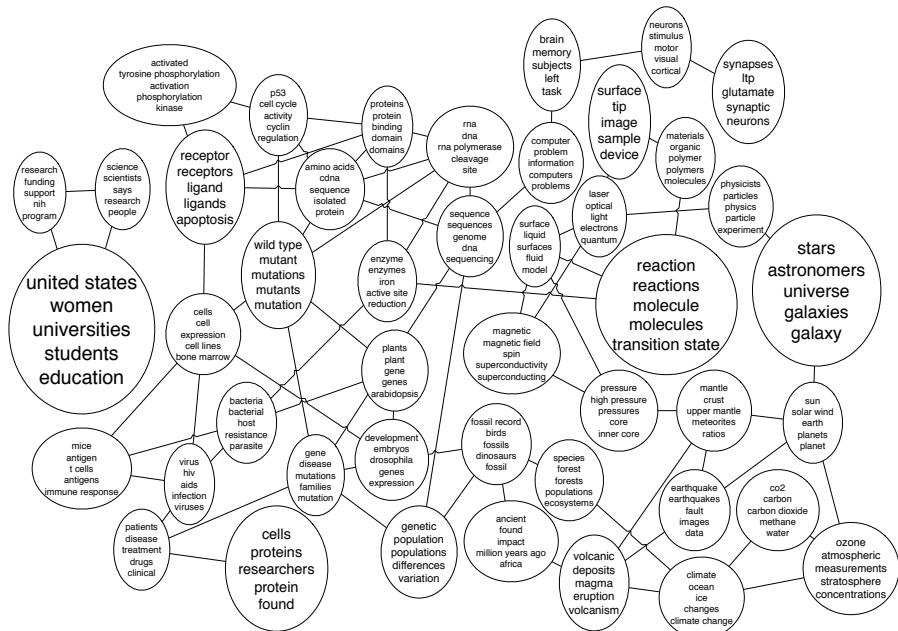
Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Latent Dirichlet Allocation

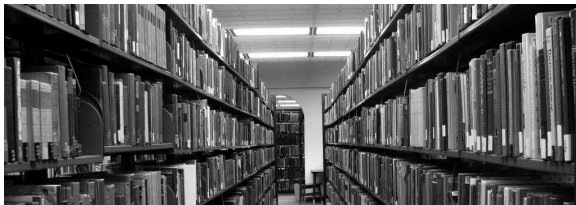


- ▶ **Data:** The OCR'ed collection of *Science* from 1990–2000
 - 17K documents
 - 11M words
 - 20K unique terms (stop words and rare words removed)
- ▶ **Model:** 100-topic LDA model using variational inference.

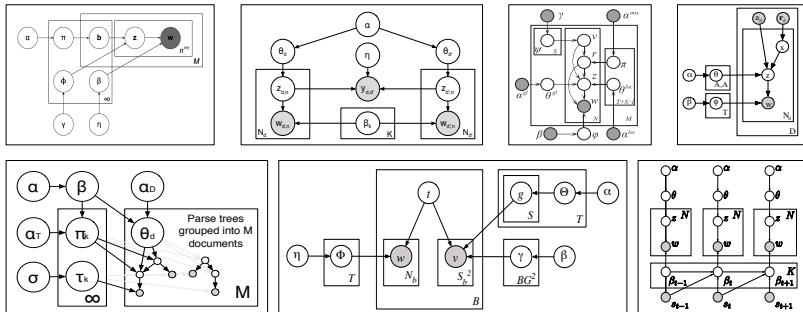




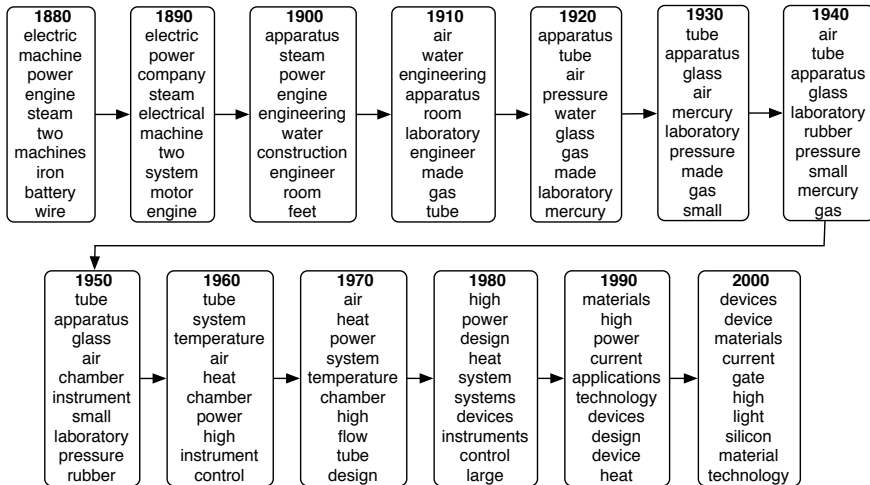
Topic Modeling



- ▶ LDA builds on decades of research about how to **derive meaning from text**.
- ▶ LDA more easily **scales to massive data** and **generalizes to new data**.
- ▶ LDA has had a big impact on many fields
 - Natural language processing
 - Computer vision
 - Recommendation systems
 - Web search
 - Computational biology and genetics



- LDA is a simple **building block** that enables many applications.
- Each model solves a different problem, fuses different kinds of data.
- Models and their algorithms easily compose.





SKY WATER TREE
MOUNTAIN PEOPLE



SCOTLAND WATER
FLOWER HILLS TREE



SKY WATER BUILDING
PEOPLE WATER



FISH WATER OCEAN
TREE CORAL



PEOPLE MARKET PATTERN
TEXTILE DISPLAY

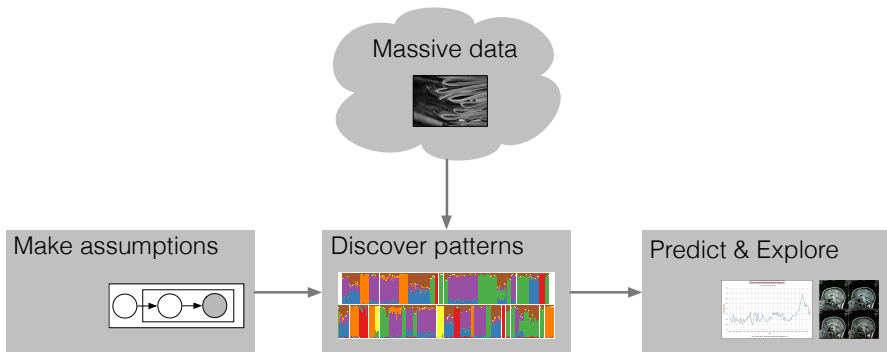


BIRDS NEST TREE
BRANCH LEAVES

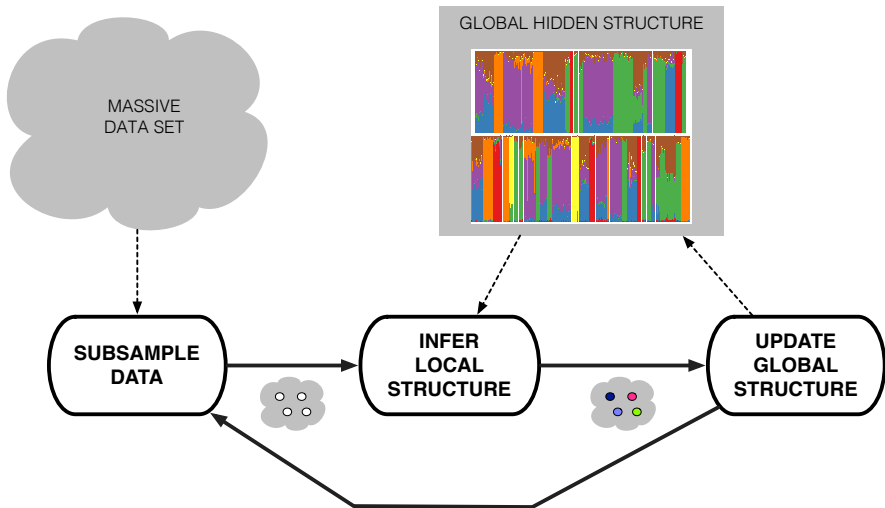


Probabilistic inference

Given a model, use an algorithm to discover the hidden patterns in the data.



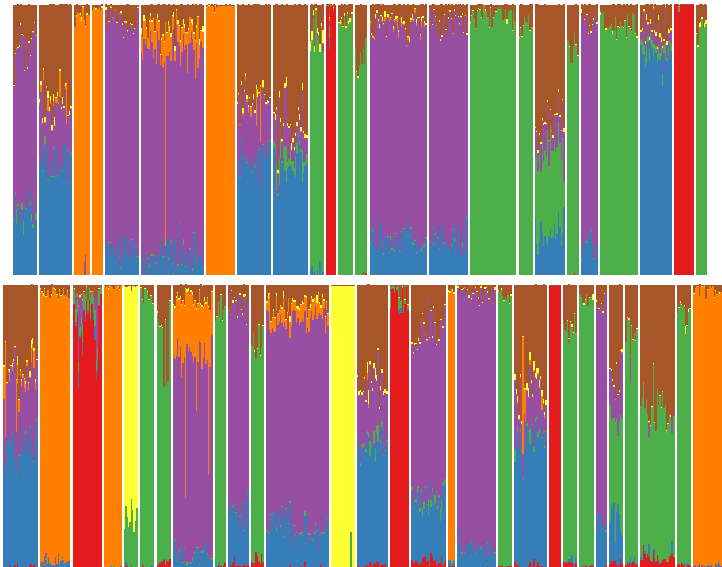
- ▶ Probabilistic inference is the main **algorithmic & statistical problem**.
- ▶ We square the modeling assumptions with the observed data.
E.g., which topics likely generated a collection of documents?
- ▶ We need **scalable** and **generic** inference.



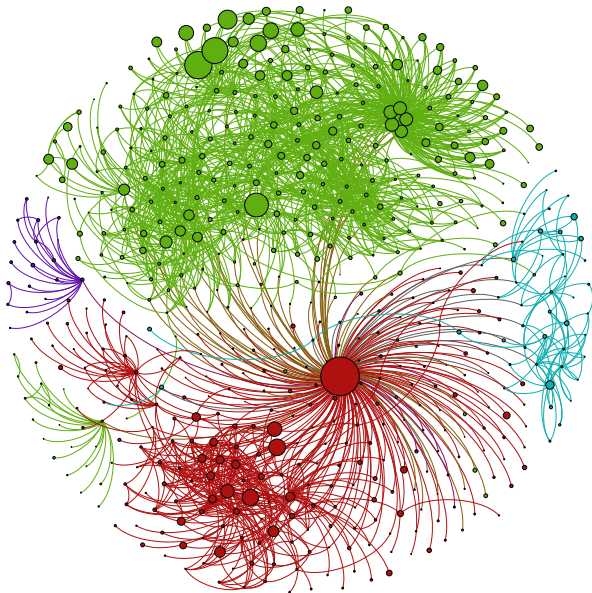
Stochastic variational inference scales to **massive data**.



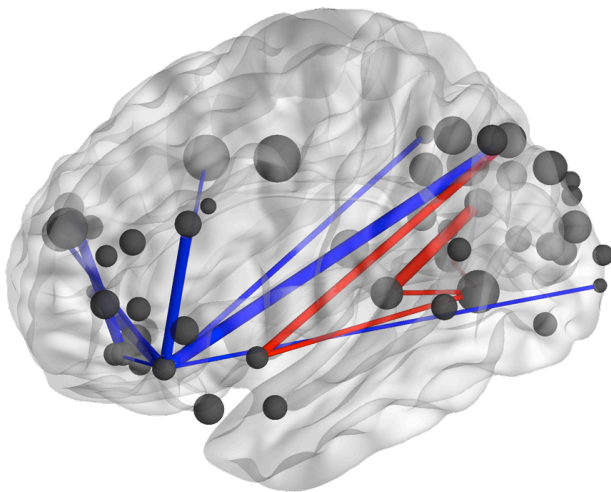
Topics found in 1.8M articles from the New York Times



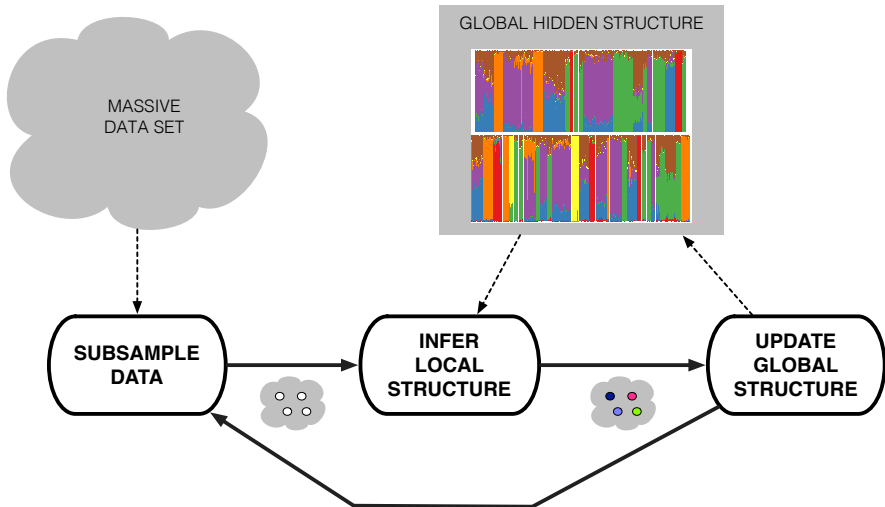
Population analysis of 2 billion genetic measurements



Communities discovered in a 3.7M node network of U.S. Patents



Neuroscience analysis of 220M fMRI measurements



- ▶ Uses **stochastic optimization** (Robbins & Monro, 1951)
- ▶ Scales up 50 years of research in Bayesian modeling
- ▶ Though these are recent results, they have been adapted to many domains

Modern probabilistic modeling:

An efficient framework for discovering useful patterns in massive data.

TOPIC
MODELING

The diagram consists of two nested ellipses. The outer ellipse is white with a black border. Inside it is a smaller, light-gray ellipse with a black border. The text 'STATISTICS', 'MACHINE LEARNING', and 'DATA SCIENCE' is centered at the bottom of the outer ellipse. Inside the gray ellipse, the text 'TOPIC MODELING' is on the left and 'PROBABILISTIC MODELING' is at the bottom. An arrow points from 'TOPIC MODELING' to a small black dot located in the upper right area of the gray ellipse.

PROBABILISTIC
MODELING

STATISTICS
MACHINE LEARNING
DATA SCIENCE

I. Assume our data come from a model with hidden patterns at work

Topics

Documents

Topic proportions and assignments

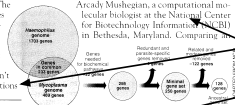


Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

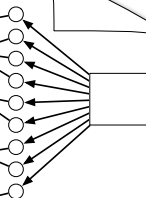
"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson, a biologist at Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a numeric game. Some particularly "core" and more genomes are repeatedly sequenced and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing the



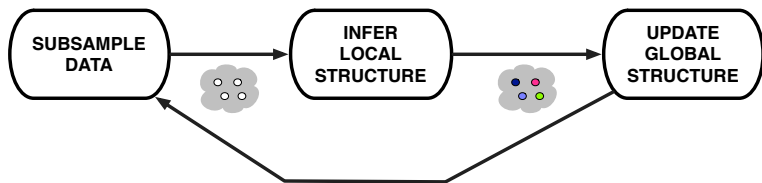
* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

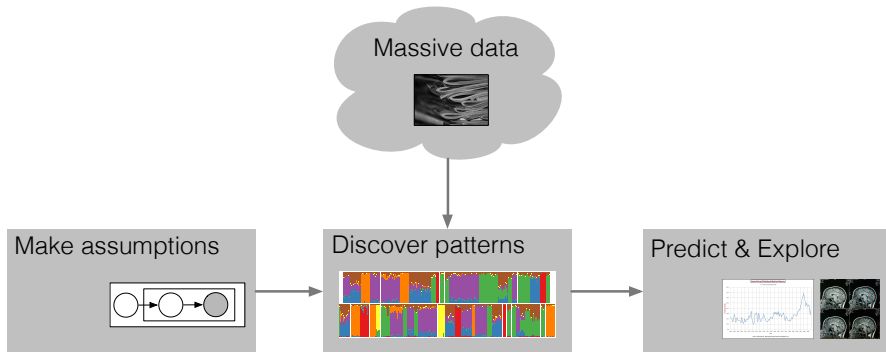
SCIENCE • VOL. 272 • 24 MAY 1996



II. Discover those patterns in the data

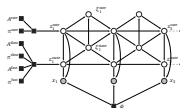
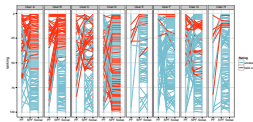
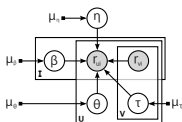
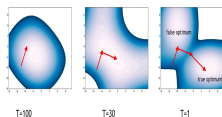


The probabilistic modeling pipeline



Our goal:

- Develop modeling into a **flexible**, **powerful** and **easy-to-use** way to solve real-world problems.



Models and applications

- ▶ genetic measurements
- ▶ hierarchical topics in a corpus
- ▶ changing preferences
- ▶ equations and text
- ▶ newsworthy events in Twitter
- ▶ news consumption in a network
- ▶ word meanings
- ▶ counselor/patient dialogs
- ▶ declassified cables from the 70s
- ▶ neural readings in a fish

Inference

- ▶ active subsampling
- ▶ averaged gradients
- ▶ annealing and inference
- ▶ stochastic optimization
- ▶ structured variational inference
- ▶ probabilistic programming



We should seek out unfamiliar summaries of observational material, and establish their useful properties... And still more novelty can come from finding, and evading, still deeper lying constraints.

(John Tukey, *The Future of Data Analysis*, 1962)